

Model	Accuracy	State updates
LSTM	0.910 ± 0.045	784.00 ± 0.00
LSTM ($p_{skip} = 0.5$)	0.893 ± 0.003	392.03 ± 0.05
Skip LSTM, $\lambda = 10^{-4}$	0.973 ± 0.002	379.38 ± 33.09
GRU	0.968 ± 0.013	784.00 ± 0.00
GRU ($p_{skip} = 0.5$)	0.912 ± 0.004	391.86 ± 0.14
Skip GRU, $\lambda = 10^{-4}$	0.976 ± 0.003	392.62 ± 26.48

TABLE I. ACCURACY AND USED SAMPLES ON THE TEST SET OF MNIST. RESULTS ARE DISPLAYED AS $mean \pm std$ OVER FOUR DIFFERENT RUNS.

e.g. normalization [3], [1], regularization [12], [7], variable computation [6], [10] or even external memory [4], [11].

Skip RNN is able to learn when to update or copy the state without explicit information about which samples are useful to solve the task at hand. However, a different operating point on the trade-off between performance and number of processed samples may be required depending on the application, e.g. one may be willing to sacrifice a few accuracy points in order to run faster on machines with low computational power, or to reduce energy impact on portable devices. The proposed model can be encouraged to perform fewer state updates through additional loss terms:

$$L_{budget} = \lambda \cdot \sum_{t=1}^T u_t \quad (6)$$

where L_{budget} is the cost associated to a single sequence, λ is the cost per sample and T is the sequence length.

III. EXPERIMENTS: SEQUENTIAL MNIST

The MNIST handwritten digits classification benchmark [9] is traditionally addressed with Convolutional Neural Networks (CNNs) that can efficiently exploit spatial dependencies through weight sharing. By flattening the 28×28 images into 784-d vectors, however, it can be reformulated as a challenging task for RNNs where long term dependencies need to be leveraged [8]. With the goal of studying the effect of skipping state updates on the learning capability of the networks, we introduce a new baseline which skips a state update with probability p_{skip} . We tune the skipping probability to obtain models that perform a similar number of state updates to the Skip RNN models.

Results in Table I show that Skip RNNs solve the task using fewer updates than their counterparts while also showing a lower variation among runs and train faster. We hypothesize that skipping updates make the Skip RNNs work on shorter subsequences, simplifying the optimization process and allowing the networks to capture long term dependencies more easily. However, the drop in performance observed in the models where the state updates are skipped randomly suggests that learning which samples to use is a key component in the performance of Skip RNN. Examples such as the ones depicted in Figure 2 show how the model learns to skip pixels that are not discriminative, such as the padding regions in the top and bottom of images, and the attended samples vary depending on the particular input being given to the network.

ACKNOWLEDGMENT

This work has been accepted as a conference paper at ICLR 2018, and we refer the reader to the full publication for extended

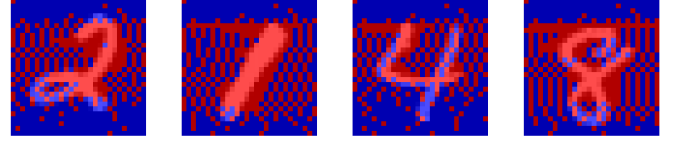


Fig. 2. Sample usage examples for the Skip LSTM with $\lambda = 10^{-4}$ on the test set of MNIST. Red pixels are used, whereas blue ones are skipped.

experiments and results [2]. This research was partially supported by the Spanish Ministry of Economy and Competitiveness and the European Regional Development Fund (ERDF) under contracts TEC2016-75976-R and TIN2015-65316-P, by the BSC-CNS Severo Ochoa program SEV-2015-0493, and grant 2014-SGR-1051 by the Catalan Government. Víctor Campos was supported by Obra Social “la Caixa” through La Caixa-Severo Ochoa International Doctoral Fellowship program. We would also like to thank the technical support team at the Barcelona Supercomputing Center.

REFERENCES

- [1] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [2] V. Campos, B. Jou, X. Giró-i Nieto, J. Torres, and S.-F. Chang. Skip rnn: Learning to skip state updates in recurrent neural networks. In *International Conference on Learning Representations*, 2018.
- [3] T. Coijmans, N. Ballas, C. Laurent, Ç. Gülçehre, and A. Courville. Recurrent batch normalization. In *ICLR*, 2017.
- [4] A. Graves, G. Wayne, and I. Danihelka. Neural Turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- [5] G. Hinton. Neural networks for machine learning. Coursera video lectures, 2012.
- [6] Y. Jernite, E. Grave, A. Joulin, and T. Mikolov. Variable computation in recurrent neural networks. In *ICLR*, 2017.
- [7] D. Krueger, T. Maharaj, J. Kramár, M. Pezeshki, N. Ballas, N. R. Ke, A. Goyal, Y. Bengio, H. Larochelle, A. Courville, et al. Zoneout: Regularizing rnns by randomly preserving hidden activations. In *ICLR*, 2017.
- [8] Q. V. Le, N. Jaitly, and G. E. Hinton. A simple way to initialize recurrent networks of rectified linear units. *arXiv preprint arXiv:1504.00941*, 2015.
- [9] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [10] D. Neil, M. Pfeiffer, and S. Liu. Phased LSTM: accelerating recurrent network training for long or event-based sequences. In *NIPS*, 2016.
- [11] J. Weston, S. Chopra, and A. Bordes. Memory networks. *arXiv preprint arXiv:1410.3916*, 2014.
- [12] W. Zaremba, I. Sutskever, and O. Vinyals. Recurrent neural network regularization. In *ICLR*, 2015.



Víctor Campos holds a BsC and a MsC degrees on Electrical Engineering from Universitat Politècnica de Catalunya. He is currently pursuing his PhD on the intersection between Deep Learning and High Performance Computing at the Barcelona Supercomputing Center, supported by Obra Social “la Caixa” through La Caixa-Severo Ochoa International Doctoral Fellowship program. His research interests focus on large scale machine learning.